

Durham Research Online

Deposited in DRO:

12 February 2020

Version of attached file:

Accepted Version

Peer-review status of attached file:

Peer-reviewed

Citation for published item:

Li, Zhonghua and Wang, Pengcheng and You, Chunyuan and Yu, Jiwen and Zhang, Xiangnan and Yan, Feilin and Ye, Zhengxiu and Shen, Chao and Li, Baoqi and Guo, Kai and Liu, Nian and Thyssen, Gregory N. and Fang, David D. and Lindsey, Keith and Zhang, Xianlong and Wang, Maojun and Tu, Lili (2020) 'Combined GWAS and eQTL analysis uncovers a genetic regulatory network orchestrating the initiation of secondary cell wall development in cotton.', *New phytologist.*, 226 (6). pp. 1738-1752.

Further information on publisher's website:

<https://doi.org/10.1111/nph.16468>

Publisher's copyright statement:

This is the accepted version of the following article: Li, Zhonghua, Wang, Pengcheng, You, Chunyuan, Yu, Jiwen, Zhang, Xiangnan, Yan, Feilin, Ye, Zhengxiu, Shen, Chao, Li, Baoqi, Guo, Kai, Liu, Nian, Thyssen, Gregory N., Fang, David D., Lindsey, Keith, Zhang, Xianlong, Wang, Maojun Tu, Lili (2020). Combined GWAS and eQTL analysis uncovers a genetic regulatory network orchestrating the initiation of secondary cell wall development in cotton. *New Phytologist* 226(6): 1738-1752 which has been published in final form at <https://doi.org/10.1111/nph.16468>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for self-archiving.

Additional information:

Use policy

The full-text may be used and/or reproduced, and given to third parties in any format or medium, without prior permission or charge, for personal research or study, educational, or not-for-profit purposes provided that:

- a full bibliographic reference is made to the original source
- a [link](#) is made to the metadata record in DRO
- the full-text is not changed in any way

The full-text must not be sold in any format or medium without the formal permission of the copyright holders.

Please consult the [full DRO policy](#) for further details.

PROF. KEITH LINDSEY (Orcid ID : 0000-0002-7992-6804)

PROF. MAOJUN WANG (Orcid ID : 0000-0002-4791-3742)

Article type : Regular Manuscript

Combined GWAS and eQTL analysis uncovers a genetic regulatory network orchestrating the initiation of secondary cell wall development in cotton

Zhonghua Li^{1,6}, Pengcheng Wang^{1,6}, Chunyuan You², Jiwen Yu³, Xiangnan Zhang¹, Feilin Yan¹, Zhengxiu Ye¹, Chao Shen¹, Baoqi Li¹, Kai Guo¹, Nian Liu¹, Gregory N. Thyssen⁴, David D. Fang⁴, Keith Lindsey⁵, Xianlong Zhang¹, Maojun Wang^{1*} & Lili Tu^{1*}

¹National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural University, Wuhan 430070, Hubei, China.

²Cotton Research Institute, Shihezi Academy of Agriculture Science, Shihezi 832000, Xinjiang, China.

³State Key Laboratory of Cotton Biology, Cotton Institute of the Chinese Academy of Agricultural Sciences, Anyang 455000, China.

⁴Cotton Fiber Bioscience Research Unit, USDA-ARS, Southern Regional Research Center, New Orleans, LA 70124, USA.

⁵Department of Biosciences, Durham University, Durham, DH1 3LE, UK

⁶These authors contributed equally to this work.

This article has been accepted for publication and undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/NPH.16468](https://doi.org/10.1111/NPH.16468)

This article is protected by copyright. All rights reserved

*Authors for correspondence:

Maojun Wang

Email: mjwang@mail.hzau.edu.cn

Telephone number: +86-02787283955

ORCID: <https://orcid.org/0000-0002-4791-3742>

Lili Tu

Email: lilitu@mail.hzau.edu.cn

Telephone number: +86-02787283955

ORCID: <https://orcid.org/0000-0003-4294-1928>

Received: *2 September 2019*

Accepted: *28 January 2020*

Summary

- The cotton fiber serves as a valuable experimental system to study cell wall synthesis in plants, but our understanding of the genetic regulation of this process during fiber development remains limited.
- We performed a genome-wide association study (GWAS) and identified 28 genetic loci associated with fiber quality in allotetraploid cotton. To investigate the regulatory roles of these loci, we sequenced fiber transcriptomes of 251 cotton accessions and identified 15,330 expression quantitative trait loci (eQTL).
- Analysis of local eQTL and GWAS data prioritized 13 likely causal genes for differential fiber quality in a transcriptome-wide association study (TWAS). Characterization of distal eQTL revealed unequal genetic regulation patterns between two subgenomes, highlighted by an eQTL hotspot (Hot216) that establishes a genome-wide genetic network regulating the expression of 962 genes. The primary regulatory role of Hot216, and specifically the gene encoding a KIP-related protein, was found to be the transcriptional regulation of genes responsible for cell wall synthesis, which contributes to fiber length by modulating the developmental transition from rapid cell elongation to secondary cell wall synthesis.
- This study uncovers the genetic regulation of fiber-cell development and reveals the molecular basis of the temporal modulation of secondary cell wall synthesis during plant cell elongation.

Key words: cotton fiber, cell wall synthesis, cell elongation, eQTL, genetic regulation

Introduction

Phenotypic diversity in a species is usually determined by a variety of genetic variations which can be investigated by genetic approaches such as genome-wide association studies (GWAS). The molecular roles of causal genetic variations in controlling phenotypes can be explored by investigating how they modulate differential expression of genes (Cookson *et al.*, 2009). At a population level, gene expression variation can be measured quantitatively and mapped to a genomic locus (Cheung and Spielman, 2009). With the application of high-throughput technologies for gene expression profiling, expression quantitative trait locus (eQTL) mapping has emerged as a potentially powerful new approach to investigate the genetic architecture of expression variation in a variety of organisms, and to provide molecular links between genetic variation and phenotypic diversity (Aguet *et al.*, 2017; Hormozdiari *et al.*, 2018; Lappalainen *et al.*, 2013; Zhu *et al.*, 2016).

Over the past decade, eQTL mapping has been used to investigate the genetic architecture of regulatory variation in gene expression in model plants and major crops (Wang *et al.*, 2010; X. Wang *et al.*, 2018; West *et al.*, 2007; L. Zhang *et al.*, 2017; Zhang *et al.*, 2011). A comprehensive eQTL analysis allows the construction of regulatory networks (Fu *et al.*, 2013; Keurentjes *et al.*, 2007; Wang *et al.*, 2014). Furthermore, characterization of eQTL provides an approach to address additional biological questions. In *Arabidopsis*, eQTL mapping has been conducted to explore the relationship between genes responding to the environment and genes regulated by genomic variation affecting their expression (Lowry *et al.*, 2013). In maize, an eQTL analysis has shown a critical role for non-coding genomic sequences in regulating expression variation, and a multi-tissue eQTL analysis revealed the contribution of rare genetic variants to expression extremes (Kremling *et al.*, 2018; Liu *et al.*, 2017). Recently, an integrated analysis of eQTL and metabolic QTL revealed the metabolic breeding history of tomato (Zhu *et al.*, 2018).

Cotton (*Gossypium* spp.) is globally cultivated for the utilization of natural renewable fibers in textiles, which dictates that the major goal of cotton breeding is to cultivate varieties producing fibers of superior quality. Genome-enabled breeding serves as an efficient approach to fulfill this goal, but requires an understanding of the genetic basis underpinning fiber quality-related traits and how genetic variations influence fiber development. The cotton fiber undergoes a staged cellular development to form a mature lint fiber after the ovule epidermal cell initiates an elongation process.

This has been used as an excellent experimental system for studying polarized cell growth and cell wall biosynthesis in plants (Haigler *et al.*, 2012). Numerous genetic mapping efforts have identified many loci that contribute to fiber quality-related traits (Fang *et al.*, 2017; Huang *et al.*, 2017; Liu *et al.*, 2018; Ma *et al.*, 2018; Thyssen *et al.*, 2019; Wang *et al.*, 2017), and functional analysis has confirmed important roles for candidate genes (Shan *et al.*, 2014; Tan *et al.*, 2013). It has also been found that most characterized genes play a role in more than one biological process, such as redox homeostasis, sugar transport and cell wall metabolism (Guo *et al.*, 2016; Han *et al.*, 2013; Li *et al.*, 2016; Z. Zhang *et al.*, 2017). However, some intriguing questions remain poorly understood, such as how those genes work together to regulate the complex process of fiber development and how they determine the transition from rapid cell elongation to secondary cell wall biosynthesis.

Here, we present an analysis of fiber transcriptomes from a natural *G. hirsutum* population with 251 accessions. This allowed us to identify 15,330 eQTL associated with 9,282 genes, and define the eQTL network for the elongating fiber. An eQTL hotspot has been characterized for regulating the expression of genes responsible for cell wall biosynthesis, which controls the trait of fiber length. This study opens the way for linking regulatory variants to gene transcription in cotton fiber development and provides a useful reference for accelerating the genetic improvement of fiber quality.

Materials and Methods

Plant Materials

In a previous study, we constructed a genomic variation map by sequencing a natural population of Upland cotton (*Gossypium hirsutum*) accessions (Wang *et al.*, 2017). To understand the regulatory mechanisms of fiber development, a total of 251 Upland cotton accessions were cultivated in the field at Shihezi (44°20'15''N, 86°3'28''E), Xinjiang, China in 2016. Shihezi has a suitable climate for cotton growth and is one of major locations for cotton cultivation in China. Cotton bolls were tagged on the day of flowering as 0 DPA (day post anthesis) samples. Three replicate populations of each accession were planted in three separate plots, with each replicate population arranged randomly in each plot. For each accession, fibers at 15 DPA from at least 10 cotton bolls of different plants were

collected and bulked together, and were immediately frozen in liquid nitrogen. Fiber samples of individual replicate populations were used for RNA extraction.

RNA extraction and sequencing

For each accession, fiber samples were isolated for RNA extraction. Total RNA was extracted using a Spectrum™ Plant Total RNA Kit (Sigma, STRN250). These RNA samples were used to construct sequencing libraries using Illumina TruSeq RNA Library Preparation Kit (Illumina, San Diego, CA, USA). Each library was sequenced on an Illumina HiSeq 4000 platform (pair-end 150 bp).

RNA-Seq data mapping and analysis

The raw RNA sequencing data were subject to adapter filtering and low-quality bases were removed using Trimmomatic (version 0.32) (Bolger *et al.*, 2014). The clean data were mapped to the reference genome sequence of *G. hirsutum* using HiSAT2 software (version 2.1.0) with default settings (Kim *et al.*, 2015; Wang *et al.*, 2019). The mapping reads were sorted to filter those reads representing PCR duplicates. Sequencing reads with mapping quality of less than 25 were filtered using Samtools (version 0.1.19) (Li *et al.*, 2009). Sequencing reads that were mapped to multiple genomic loci were filtered from the mapping files using customized Perl scripts. The remaining reads were used to calculate the expression levels of genes (fragments per kilobase of transcript per million mapped reads, FPKM) using Stringtie software (version 1.3.4) with default settings (Pertea *et al.*, 2015).

Identification of genomic SNPs

In this study, whole genome re-sequencing data of 251 accessions which were used for RNA-Seq analysis were mapped to the new reference genome sequence of *G. hirsutum* using BWA software (version 0.7.10-r789) using the mem method (-M -k 25) (Li and Durbin, 2009; Wang *et al.*, 2019). The unique mapping reads were parsed to identify SNPs with Samtools and GATK (version 3.1.1) as described previously (Li *et al.*, 2009; McKenna *et al.*, 2010; Wang *et al.*, 2017). Only those SNPs that were supported by both software tools were retained. The remaining SNPs were subject to a filtering process based on the mapping depth (>6×) and quality score (>100). VCFtools (version 0.1.15) was used to merge SNPs from each accession, of which those with a minor allele frequency (MAF) of less

than 5% were filtered (Danecek *et al.*, 2011). The ANNOVAR software was used to annotate the remaining SNPs, including those representing genic variations (synonymous, non-synonymous, intronic, un-translated regions, upstream 1Kb sequence of the transcription start site and downstream 1 Kb sequence of transcription termination site) and intergenic variations (Hakonarson *et al.*, 2010).

GWAS on fiber quality-related traits

In a previous study, data on fiber quality-related traits for a natural population, including fiber length, fiber strength, elongation, uniformity and micronaire value, were collected over three years (Wang *et al.*, 2017). In the current study, we collected data for these traits in 2016 and 2017, and conduct a GWAS analysis. Compared with our previous study (Wang *et al.*, 2017), the GWAS analysis here leveraged new phenotypic data measured in the same year as when fiber samples were collected for RNA-Seq. GWAS was performed using genomic variation data (2,372,767 SNPs) from those 251 accessions for eQTL analysis. Association analysis was carried out using TASSEL (version 5.0) and FaST-LMM (version 2.02) programs (Kroon *et al.*, 2007; Lippert *et al.*, 2011). The population structure was calculated using the structure program and the kinship was derived from all SNPs (Falush *et al.*, 2003). The cutoff for determining significant associations was $P < 4.2 \times 10^{-7}$ ($1/n$, where n represents the total number of genomic SNPs). To help identify putative candidate genes for GWAS loci, RNA-Seq data for different fiber developmental stages were used to examine gene expression levels (Zhang *et al.*, 2015).

Identification of expression QTL (eQTL)

To identify eQTL for genes in fiber development, an analysis of gene expression levels was performed to discard those with expression levels less than 0.1 (FPKM < 0.1) in more than 95% of samples. This allowed us to filter 30,336 of 70,199 genes in the reference genome. The remaining genes were subject to a filtering process which requires an expression change of at least two fold between two samples representing the 5th percentile and the 95th percentile of sorted expression levels, respectively. In this step, 4,098 genes were filtered. The expression levels of remaining genes were normalized using a normal quantile transformation. In total, we obtained 35,765 genes with transformed expression levels which were regarded as expression traits for GWAS analysis. The

method of eQTL identification was similar to a previous study (Fu *et al.*, 2013), with minor modifications. In the current study, the merged genomic SNPs (MAF >0.05) were used to perform GWAS for each gene using the FaST-LMM program (Lippert *et al.*, 2011). The cutoff for determining significant associations was the same as that for GWAS. The significant SNPs for each trait (gene) were grouped into clusters with a maximum distance of 10 Kb between two consecutive SNPs, and only those clusters with more than 3 significant SNPs were retained as putative eQTL. Each eQTL was represented by the most significant SNP of all (lead SNP). The putative eQTL which were identified in a LD region were further filtered as false positive associations. In addition, to identify eQTL hotspots, the hot_scan program was run for all eQTL in each chromosome (-m 5000 -s 0.05) (Holanda *et al.*, 2014).

TWAS

The idea of transcriptome-wide association studies (TWAS) was described in a previous study (Gusev *et al.*, 2016), which provided an approach for identifying significant *cis* genetic correlation between expression and phenotypic traits. This method used a reference population with available gene expression and genetic variation data to impute the *cis* genetic component of expression into another set of phenotyped individuals for which genetic variation has been identified. The fusion program was used to conduct TWAS in this study (<http://gusevlab.org/projects/fusion/>). Since no measured population expression data in cotton has been generated previously, the normalized expression data of 3,690 genes with local eQTL and SNP data from upstream 500 Kb to downstream 500 Kb of these genes were extracted to compute expression weights, which represent the pre-modeled relationship between SNPs and expression levels of genes (Gusev *et al.*, 2016). In this analysis, heritability was calculated using GCTA (version 1.26.0) and expression weights were computed with the option --models top1, blup, lasso, enet (Yang *et al.*, 2011). The genome-wide GWAS SNP data were used as the reference data. The GWAS statistics Z-score was computed based on GWAS p-values and effect size for each respective fiber quality-related trait. The TWAS was performed on each chromosome with significant GWAS associations for traits. The TWAS p-values were corrected with the Bonferroni method in R (version 3.2.3).

Network construction

In the eQTL hotspot data, we found that some genes were regulated by several eQTL. To systematically understand the regulation associations, we leveraged a genetic network involving genes (traits) and eQTL data to visualize the relationship between eGenes and eQTL. eQTL hotspots were highlighted in the genetic network. This network included the relationship between eQTL and their regulated genes, and between eQTL hotspots and regulated genes. In this analysis, those genes which included the lead SNPs of eQTL were regarded as causal regulators. The construction of the network was performed using Cytoscape (version 3.6.1) with the edge-weighted force directed layout method (Shannon *et al.*, 2003).

Differential expression and gene set enrichment analyses

The differential expression analysis of genes in cotton accessions with extreme long fiber and extreme short fiber was performed using the DESeq package in R with a significance of false discovery rate ($FDR \leq 0.05$) (Anders and Huber, 2010). The gene ontology (GO) enrichment analysis of genes in the eQTL network or exhibiting differential expression was performed using the AgriGO webserver (version 2.0) with genes showing expression levels in fiber development as a reference (Tian *et al.*, 2017). GO terms with a FDR threshold of 0.05 were considered as significant terms.

Data availability

The raw RNA-Seq data generated in this study are available in the BioProject database under accession number PRJNA433615.

Results

QTL Identification for fiber quality-related traits using GWAS

In a previous study, we generated genomic re-sequencing data for a natural population, which provided a resource to identify favorable variants associated with fiber quality-related traits (Wang *et al.*, 2017). In the present study, another 15 accessions were sequenced and integrated with previous

data to generate an individual set of 251 accessions. We aligned sequence data of the 251 accessions against our recently published reference genome sequence of *G. hirsutum* acc. Texas Marker-1 (TM-1) (Wang *et al.*, 2019), and identified 2,372,767 SNPs with a minor allele frequency of > 0.05. We then used the multiple phenotypic data for fiber quality, including those measured under the same environmental conditions used for sample collection for transcriptome sequencing in 2016 (**Table S1**), to perform a GWAS analysis of fiber quality-related traits, including fiber length (FL), strength (FS), elongation (FE), uniformity (FU) and micronaire value (MV). A total of 28 associations were identified for these traits, including 10, 6, 8, and 4 loci for length, strength, elongation and uniformity, respectively (**Fig. 1; Table S2, S3**). No significant candidate loci were identified for micronaire value. Among these associations, 16 were previously uncharacterized (**Table 1**).

Transcriptome sequencing and eQTL mapping

GWAS analysis identified candidate loci associated with fiber quality, but how these loci contribute to different fiber properties between accessions remains largely unknown. We sought to find an answer to this question by exploring expression variation of genes that may be influenced by those loci across a population. We sequenced fiber transcriptomes for 251 *G. hirsutum* accessions, with fiber samples at 15 days post anthesis (DPA), which represents a late elongation period before the developmental transition to secondary cell wall synthesis. The accessions for RNA-Seq analysis are the same as those for genome resequencing. In total, we generated 10 billion pair-end reads with an average of 40 million for each accession (**Table S4**). These RNA-Seq data were mapped to the TM-1 reference genome to quantify the expression levels of genes. In this analysis, we found 39,863 genes with expression levels at this stage of fiber development. A total of 35,765 genes that exhibited expression variation, representing 50.9% of the annotated genes (70,199) in TM-1 genome, were used for subsequent analysis.

We performed eQTL mapping using the same genomic SNP dataset as that for GWAS analysis of fiber quality-related traits. In total, we identified 15,330 eQTL associated with the expression of 9,282 genes (eGenes regulated by eQTL). These eQTL and eGenes were from all 26 chromosomes (**Fig. 2a**). We found that the associations of eQTL and eGenes located on the same chromosome had a higher significance than those located on different chromosomes (two-sided Wilcoxon rank sum test, *P*-value

$< 2.2 \times 10^{-16}$; **Fig. S1**). For the inter-chromosomal associations, those that occurred in homoeologous chromosomes between A- and D-subgenomes (691 of 15330) were prevalent (Fisher exact test, P -value $< 2.2 \times 10^{-16}$; **Fig. 2a** shown by two red arrows). One intriguing observation is that the associations between chromosome D11 and other chromosomes (962 of 15330) were enriched (Fisher exact test, P -value $< 2.2 \times 10^{-16}$; **Fig. 2a** shown by an orange arrow).

Based on the distance between eQTL and eGenes, we categorized all eQTL into 5,370 local eQTL (< 1 Mb) and 9,960 distal eQTL (> 1 Mb or in different chromosomes). We found that local eQTL had a larger effect on expression variation than did distal eQTL (two-sided Wilcoxon rank sum test, P -value $< 2.2 \times 10^{-16}$; **Fig. 2b**), which agrees with similar findings in other organisms (Wang *et al.*, 2010; X. Wang *et al.*, 2018; L. Zhang *et al.*, 2017; Zhang *et al.*, 2011). It was also found that the distance between local eQTL and eGenes shows a preferential distribution at ca. 5 Kb (**Fig. 2c**). Of the total eQTL, local eQTL account for 35%, while distal eQTL account for 65%, of which 15% occurred on the same chromosome (distal_intraChr) and the other 50% were found on different chromosomes (distal_interChr). The distal_interChr eQTL were further divided into three groups, i.e., within the A-subgenome (At-At; 9%), within the D-subgenome (Dt-Dt; 15%), and between the A- and D-subgenomes (At-Dt; 26%) (**Fig. 2d**). Of the eGenes, we found that 5,027 were regulated by local eQTL and 6,220 by distal eQTL, and the majority of eGenes (6,049, 65.1%) were regulated by only one eQTL (**Fig. 2e**).

Global roles of eQTL in the genetics of fiber quality-related traits

Recent studies demonstrate that transcriptome-wide association studies (TWAS) represent a powerful approach for prioritizing causal genes for GWAS loci, using the information of *cis*-eQTL (Gusev *et al.*, 2016). TWAS establishes a connection between gene expression, one kind of molecular phenotype, and other physical phenotypes in organisms (Gusev *et al.*, 2018). Here, on the basis of characterizing local eQTL, we attempted to prioritize likely causal genes for fiber quality-related GWAS loci using TWAS.

We modeled 3,690 expression weights using SNPs from genomic region 500 Kb both upstream and downstream of each eGene that was transcriptionally regulated by local eQTL. These expression weights were leveraged to perform a TWAS on each chromosome with GWAS loci. The TWAS

identified 10 transcriptome-wide-significant associations, including four for fiber length, three for fiber elongation and three for fiber uniformity (Bonferroni corrected P -value <0.05 ; **Table 2; Fig. S2**). In these associations, we found that a MYB transcription factor gene (*Ghir_D05G027990*) was prioritized to be a candidate gene for a pleiotropic locus associated with fiber length and fiber uniformity. We also noted that the TWAS prioritized two genes (*Ghir_D04G017090* and *Ghir_D04G016300*) for a fiber elongation-related GWAS association on chromosome D04, and two genes (*Ghir_D12G015450* and *Ghir_D12G015670*) for a fiber uniformity-related GWAS association on chromosome D12. Even though the TWAS failed to identify gene-trait associations for all GWAS loci, this approach represents an effective way to prioritize likely causal genes in the genetics of fiber quality-related traits.

Unequal subgenome transcription regulation as revealed by eQTL analysis

The hybridization event between two diploid species that occurred 1-2 million years ago resulted in the appearance of allotetraploid *G. hirsutum*, which significantly increased regulatory complexity of gene transcription with an additional feature of inter-subgenomic transcription regulation relative to diploid ancestors (Adams *et al.*, 2003). However, few studies have been carried out to estimate the effect of inter-subgenomic regulation on gene transcription in cotton (M. Wang *et al.*, 2018). The mapping of distal eQTL provides an avenue to investigate inter-subgenomic genetic regulation for gene transcription in fibers. In this study, we found that 57.1% of eQTL in the A-subgenome, including 52.6% for inter-subgenomic regulation and 4.5% for both inter- and intra-subgenome regulations, are associated with the transcription of genes in the D-subgenome. An analysis of eQTL in the D-subgenome showed 58.2% (46.5% and 11.7%) of eQTL are responsible for transcriptional regulation of genes in the A-subgenome (**Fig. 3a**).

In terms of the number of genes, we found that 44.3% of the eGenes identified in the A-subgenome are regulated by eQTL in the D-subgenome, including 29.9% regulated by inter-genomic eQTL and 14.4% regulated by both inter- and intra-subgenomic eQTL. However, an analysis of eGenes in the D-subgenome shows that only 23.4% of eGenes (9.9% and 13.5%) have eQTL regulation in the A-subgenome (**Fig. 3b**). This difference indicates that a larger number of genes in the A-subgenome are transcriptionally regulated by the D-subgenome at this fiber

developmental stage (Pearson's Chi-squared test, P -value $<2.2 \times 10^{-16}$), highlighting unequal transcriptional regulation patterns between the two subgenomes. A gene ontology (GO) enrichment analysis shows that eGenes in the A-subgenome with inter-subgenomic eQTL regulation are involved in biological processes linked to fiber development, such as REDOX homeostasis, plant cell wall organization or biogenesis, and cell tip growth (**Fig. 3c; Table S5**). In comparison, eGenes in the D-subgenome that are regulated by the A-subgenome are enriched in fundamental processes such as regulation of DNA metabolism and neutral amino acid transport.

To specifically explore the inter-subgenomic regulation of fiber quality, we investigated the colocalization of distal eQTL and GWAS loci for fiber-quality related traits. We found that eight GWAS associations had significant SNPs that overlapped with inter-subgenomic eQTL, including two for fiber length on the D05 and D11 chromosomes, two for fiber strength on the A12 chromosome, three for fiber elongation on the A13, D04 and D05 chromosomes, and one for fiber uniformity on the D12 chromosome (**Table S6**). These colocalized eQTL regulated a total of 1,042 eGenes, of which 488 of 507 (96.2%) in the A-subgenome were regulated by eQTL in the D-subgenome and 30 of 535 (5.6%) in the D-subgenome were regulated by eQTL in the A-subgenome, which suggests a primary role for the D-subgenome in regulating genetic loci associated with fiber quality in terms of inter-subgenomic regulation. Of note is the observation that a fiber length-associated GWAS locus on chromosome D11 was colocalized with eQTL regulating 479 eGenes in the A-subgenome. These data suggest that further characterization of the functional roles of these GWAS loci should be integrated with analysis of their inter-subgenomic regulatory roles in gene transcription.

Characterization of eQTL hotspots and a regulatory network for fiber length

As shown in previous studies, some eQTL are located in a genomic region which affects the expression of many genes, and such a region is known as an eQTL hotspot (Liu *et al.*, 2017; X. Wang *et al.*, 2018). We investigated whether there is a similar phenomenon in cotton fiber development, and a total of 243 hotspots were identified (**Table S7**), regulating the expression of 3,820 genes (41.1% of all identified eGenes). 125 of these eQTL hotspots were identified in the A-subgenome and 118 in the D-subgenome. For each hotspot, the number of eGenes varies from 3 to 962 (**Table S7**). We predicted

possible key regulators for these hotspots affecting the expression of downstream eGenes by integrating genetic position information for both eQTL and gene annotations.

The identification of eQTL hotspots allows the characterization of the complex regulatory relationship between eQTL and eGenes. Here, we constructed an eQTL network consisting of 44 eQTL hotspots, 220 eQTL and 1,896 eGenes (**Fig. 4a**). In this network, the eQTL hotspot 216 (Hot216) on chromosome D11 (24.43–24.62 Mb) was highlighted, and was found to regulate the expression of 962 genes, including 479 in the A-subgenome and 483 in the D-subgenome (**Table S8**). This explains why the genome-wide analysis shows that expression variation of a large number of genes is associated with a genomic region in chromosome D11 (**Fig. 2a**). A detailed analysis of those genes shows that the expression of 293 genes is also regulated by other eQTL or hotspots, and the remaining 669 genes are only regulated by Hot216 (**Fig. 4b**).

To further characterize the biological role of Hot216, we performed a GO enrichment analysis of the 962 genes. It was observed that these genes show an enrichment in microtubule cytoskeleton for the category of cellular component (CC); protein kinase and cellulase activity for the category of molecular function (MF); and cell wall organization or biogenesis for the category of biological process (BP; **Fig. 4c**; **Table S9**). Given these relevant predicted functions, the large number of eGenes regulated by the Hot216 indicates that it represents a powerful eQTL hotspot for the transcriptional regulation of fiber development.

The effect of genomic and transcriptional variability on fiber length

To characterize the details of the 962 genes regulated by the Hot216, we performed an expression analysis across the 251 accessions. It was found that the expression patterns of these genes cluster into two groups (group-I and group-II; **Fig. 5a**). Group-I includes 287 genes which were highly expressed in 91 accessions and group-II includes the other 675 genes which were highly expressed in the other 160 accessions. GO analysis showed that genes in group-I are enriched in biological processes of plant-type cell wall organization or biogenesis (GO:0071669) and cell wall biogenesis (GO:0042546), while genes in group-II are enriched in processes including activation of MAPKK activity (GO:0000186), developmental growth (GO:0048589) and cell tip growth (GO:0009932). We

hypothesized that the significant GO terms for group-II represent processes which may promote fiber cell elongation and development.

We then compared the fiber quality-related traits of accessions between cluster-I and cluster-II. It was found that fiber length in cluster-I is significantly shorter than that in cluster-II (two-sided Wilcoxon rank sum test, P -value $< 7.4 \times 10^{-5}$), and the other three traits show no significant difference between the two clusters (**Fig. S3**). This suggests that the analysis of gene expression patterns categorized the accessions into two clusters with significant differences of fiber length. Accessions in cluster-I exhibit a short fiber phenotype with high expression of genes responsible for cell wall synthesis, while accessions in the cluster-II produce longer fiber accompanied by high expression of genes associated with cell growth. This result is supported by previous functional characterization of genes amongst these 962 eGenes (**Fig. 5b**). For example this includes three genes encoding two fasciclin-like arabinogalactan proteins (FLA7: Ghir_A08G005490 and FLA11: Ghir_D11G035910) and a trichome birefringence-like protein (TBL3: Ghir_D13G010200; Ghir_A04G010010), with established roles in secondary-wall cellulose synthesis (Bischoff *et al.*, 2010; MacMillan *et al.*, 2010). Two genes encoding an irregular xylem-related protein (IRX9: Ghir_A09G016060) and a xyloglucan endotransglucosylase/hydrolase (XTH30: Ghir_A08G016210) have been shown to be involved in xylan and xyloglucan metabolic processes respectively (Bourquin *et al.*, 2002; Lee *et al.*, 2007). We also highlighted three MYB transcriptional factors (MYB46: Ghir_A13G022890, MYB61: Ghir_A07G014020 and MYB103: Ghir_A08G012250; Ghir_D08G012890) which are involved in the positive regulation of secondary cell wall biogenesis (Kim *et al.*, 2013; Taylor-Teeples *et al.*, 2015).

To understand whether gene expression patterns of the Hot216-guided network have a genetic association with fiber length, we overlapped this hotspot region with GWAS signals. It is found that the Hot216 has the same genomic location as a significant GWAS locus for fiber length on chromosome D11 (24.44–24.62 Mb; **Fig. 1**; **Fig. 5c**). Integration of eQTL and GWAS data led us to identify a candidate gene encoding KIP-related protein 6 (KRP6) which contains a non-synonymous mutation in the first exon (G to T transition). Interestingly, the expression of *KRP6* was also regulated by the same genetic region which might act in *cis*-eQTL regulation (**Fig. 5d**). Genetic modification of KIP-related protein in *Arabidopsis* has been found to affect the expression of genes involved in plant cell wall organization and heterochromatin modification, and regulated cell elongation (Jégu *et al.*,

2013). Here, we found that two different genotypes of *KRP6* correspond to differential fiber length (**Fig. 5e**). Moreover, *KRP6* exhibits differential expression levels between accessions with two different genotypes (**Fig. 5f**). In addition, many eGenes in Hot216, known to be involved in secondary cell wall organization or biogenesis, such as *FLA7/11* and *MYB46/103*, exhibit differential expression levels between accessions with different genotypes of *KRP6* (**Fig. 5f**). These results indicate that *KRP6* is a candidate gene for differential fiber cell length via transcriptional regulation of a large number of genes in fiber development.

To further support the view that differential expression of genes responsible for cell wall synthesis may contribute to varied fiber length, we made a comparative transcriptome analysis of 60 accessions with extreme fiber length. 30 accessions with short fibers were found to exhibit up-regulated expression of 1,163 genes and down-regulated expression of 507 genes compared with 30 accessions with long fibers (**Fig. S4a**). A GO analysis showed that these up-regulated genes were enriched in plant-type cell wall organization or biogenesis, xylan biosynthetic process and microtubule depolymerization, while the down-regulated genes were enriched in categories designated response to stimulus and inorganic anion transport (**Fig. S4b**). We found that genes encoding cellulose synthase A catalytic subunit (*CesA*) 4, 7 and 8, which are required for secondary cell wall biosynthesis (Taylor-Teeple *et al.*, 2015), are highly expressed in accessions with short fibers (**Fig. 6**). Even though these *CesAs* were not found to be directly regulated by the Hot216, they act downstream of secondary wall cellulose synthesis and may be transcriptionally regulated by other genes with direct regulation from the Hot216 such as *MYB46* (Kim *et al.*, 2013; Taylor-Teeple *et al.*, 2015).

We conclude from these data that genetic variations at the GWAS locus on chromosome D11, specifically the candidate gene *KRP6*, can induce the expression of genes responsible for cell wall synthesis, and contributes to the early biosynthesis of secondary cell wall that leads to the formation of a shorter cotton fiber cell.

Discussion

eQTL analysis links regulatory variants to gene transcription

With the recent advances in sequencing the cotton genome, a very large number of genetic variants have been identified in different accessions, but only a few loci with variants have been found to be associated with agronomic traits by GWAS (Fang *et al.*, 2017; Huang *et al.*, 2017; Liu *et al.*, 2018; Ma *et al.*, 2018; Thyssen *et al.*, 2019; Wang *et al.*, 2017). It is believed that many other genetic variants may have a regulatory role in gene expression, but their regulatory targets remain undetermined. In a previous study, we integrated DNase I digestion followed by sequencing (DNase-Seq) and high-throughput chromosome conformation capture (Hi-C) data to annotate regulatory variants in *cis*-regulatory elements of promoters and distal enhancers of genes (Wang *et al.*, 2017). Here, we demonstrate that eQTL mapping represents another high-throughput approach to link regulatory variants to gene expression in fiber development. Our mapping of 15,330 eQTL represents the discovery of regulatory variants only at the transition from cell elongation to secondary cell wall development. Nevertheless, it would be expected that more regulatory variants would be identified if RNA-Seq data from different fiber developmental stages were generated. Nevertheless we demonstrate the practical application of TWAS for prioritizing likely causal genes for GWAS loci. Compared with functional analysis of lead SNPs or homology-based prediction of causal genes for GWAS loci, TWAS establishes a direct connection between gene expression and phenotype using eQTL data. This suggests that an investigation of the intermediate omics–transcriptome between variome and phenome can facilitate the understanding of the regulatory roles of genetic variants in shaping phenotypic differences.

Inter-subgenomic regulation increases regulatory complexity of gene transcription

It is known that many agronomic traits in allopolyploids, such as fiber quality in allotetraploid cotton, are regulated by the coordination of different subgenomes (Fang *et al.*, 2017; Ma *et al.*, 2018; Thyssen *et al.*, 2019; Wang *et al.*, 2017; Yuan *et al.*, 2015). Decoding the role of each subgenome in regulating desirable traits will enhance our understanding of the effects of polyploidization on plant development. In cotton, previous studies only compared the expression patterns of homoeologous genes between A- and D-subgenomes in fiber development, which uncovered subgenome expression dominance or gene expression bias (Hovav *et al.*, 2008). In this study, we used eQTL data to identify widespread inter-subgenomic regulation, and show that the D-subgenome has a relatively large regulatory effect

on the A-subgenome in terms of the number of eGenes, indicating the important role of the D-subgenome in fiber development. Even though we do not yet have similar evidence for the entire fiber development process, one of conclusions that can be made is that the D-subgenome has a large regulatory effect on the development of spinnable fiber in cultivated tetraploid cotton, because more selection signals during domestication were identified in the D-subgenome, which drove the phenotypic change from short and pigmented to long and white fibers (Wang *et al.*, 2017).

eQTL data represent genetic evidence for inter-subgenomic regulation of gene expression. Previously, we used Hi-C data to identify a number of subgenomic chromatin interactions which were proposed to play a role in the coordination of expression of homoeologous genes in tetraploid cotton (M. Wang *et al.*, 2018). It remains to be determined whether the inter-subgenomic eQTL regulation of gene transcription relies on spatio-chromatin contact between chromosomes. Interestingly, we observed a pattern of subgenome contacts using Hi-C data in leaf (**Fig. S5**), which was similar to the pattern of eQTL regulation (**Fig. 2a**). This open question should be investigated further by using high-throughput chromatin contact matrix and eQTL data from the same tissue. From an evolutionary view, both subgenomes were divergent from a common ancestor and had a high sequence similarity, so future studies should explore whether homoeologous sequences in one subgenome could regulate the expression of genes in the other subgenome. The analysis of inter-subgenomic regulation may provide a strategy for exploring the regulatory mechanism of expression novelty in polyploids and inform our understanding of phenotypic advantage.

A Hot216-guided genetic regulatory network orchestrates the initiation of secondary cell wall development in plants

After initiation, cotton fiber cells undergo a rapid elongation stage followed by a secondary cell wall synthesis stage in which cellulose is synthesized. The duration of elongation plays a vital role in determining mature fiber length, and the biosynthesis of the secondary cell wall restricts fiber elongation, as supported by metabolomic and transcriptomic data (Haigler *et al.*, 2012; Tuttle *et al.*, 2015). In this study, we integrate GWAS, TWAS, eQTL networks and transcriptome analysis to address the question of the genetic regulation of fiber length, and identify an important role for a genetic locus on chromosome D11. Even though this locus has been identified in previous studies (Ma

et al., 2018; Thyssen *et al.*, 2019), the regulatory mechanism underlying fiber length has remained largely unknown. We demonstrate that the likely causal gene *KRP6* acts as an eQTL hotspot (Hot216) to regulate the expression of 962 genes which are involved in a regulatory network. This large number of genes regulated by Hot216 provides a similar picture to results of experiments using chromatin immunoprecipitation followed by high-throughput sequencing (ChIP-Seq) of KRP5 in *Arabidopsis*, which showed an enrichment of genes involved in cell wall organization (Jégu *et al.*, 2013). This suggests a conserved regulatory role of KRPs in promoting cell elongation between cotton and *Arabidopsis*.

In plants, KRPs have been found to act as inhibitors of cyclin-dependent kinases and have a role in endoreduplication (Barrôco *et al.*, 2006; García-Ramírez *et al.*, 2017; Jégu *et al.*, 2013; Zhao *et al.*, 2017). In cotton, a previous study discovered that developing fiber cells undergo endoreduplication but not cell division (Van't Hof, 1999). The single-cell characteristic of the cotton fiber prompts us to propose that a key function of *KRP6* might be to ensure a maintained cell cycle interphase that contributes to fiber-cell elongation, whereas an abnormal function of *KRP6* or down-regulated expression would lead to an earlier transition to secondary cell wall synthesis. This possibility is supported by the identification, in this genetic regulatory network, of several up-regulated secondary cell wall synthesis genes in cotton accessions with relatively short fibers. Therefore, we propose the view that the interplay between secondary wall synthesis and endoreduplication may play a vital role in determining the rate or timing of fiber elongation and final length (Sablowski and Carnier Dornelas, 2013). Characterization of this genetic locus and the related regulatory network suggests that future improvement of fiber quality, especially for fiber length, should focus on the temporal manipulation of secondary cell wall synthesis by genetic selection of eQTL hotspots such as Hot216.

In conclusion, we provide an example of the power of integrating genotypic, gene expression and phenotypic analysis to gain insights into the genetic regulation of fiber quality-related traits, which informs our understanding of the role of secondary cell wall in modulating cell elongation in plants.

Acknowledgements

This project was financially supported by the National Natural Science Foundation of China (31830062; 31471540) to X.Z. and L.T., and was also supported by National Key Research and Development Program of China (2016YFD0100505-05) and Fundamental Research Funds for the Central Universities (2662017JC030) to L.T. Funding was also provided by the National Postdoctoral Program for Innovative Talents (BX201700094) to M.W. The authors declare no competing interests.

Author contributions

L.T., M.W. and Xianlong Zhang conceived and designed the project. C.Y. and J.Y. managed cotton cultivation in the field. P.W., Z.L., Xiangnan Zhang, F.Y., Z.Y., K.G. and N.L. collected fiber samples and performed experiments. L.T., Z.L. and P.W. performed the RNA sequencing. M.W., C.S., P.W., B.L. and Z.L. analyzed the data. G.N.T., D.D.F. and K.L. contributed to project discussion. M.W. wrote the manuscript draft. X.Z., L.T., K.L. and D.D.F. revised the manuscript. Z.L. and P.W. contributed equally to this work.

Table 1 Summary of genome-wide association study (GWAS) results for fiber quality-related traits in *Gossypium hirsutum*.

QTL	Chr	SNP number ^a	lead SNP position	Major allele	Minor allele	Minor allele%	P-value	Refs ^c
FL1	A05	4	38429540	T ^b	A	12.7%	2.0×10 ⁻⁷	
FL2	A07	8	93511292	T ^b	C	10.8%	4.1×10 ⁻⁹	Ma <i>et al.</i> , 2018
FL3	A08	2	92006893	T ^b	A	5.2%	8.8×10 ⁻¹⁰	Liu <i>et al.</i> , 2018
FL4	A11	2	45735171	A ^b	C	6.8%	2.7×10 ⁻⁷	
FL5	A13	2	103375434	T ^b	C	15.9%	2.9×10 ⁻⁷	Fang <i>et al.</i> , 2017; Liu <i>et al.</i> , 2018
FL6	D01	5	9210491	G ^b	A	8.0%	2.8×10 ⁻⁸	
FL7	D05	135	27233958	T ^b	C	5.6%	5.2×10 ⁻⁹	Ma <i>et al.</i> , 2018
FL8	D06	6	59462878	C	T ^b	21.5%	1.1×10 ⁻⁷	
FL9	D11	90	24616418	G ^b	A	17.1%	6.6×10 ⁻¹⁰	Ma <i>et al.</i> , 2018; Thyssen <i>et al.</i> , 2019
FL10	D11	3	65633365	T ^b	C	5.2%	2.4×10 ⁻⁸	
FS1	A06	1	90655795	C ^b	T	6.8%	1.8×10 ⁻⁷	
FS2	A11	7	10086019	C	T ^b	5.2%	2.8×10 ⁻⁸	Liu <i>et al.</i> , 2018
FS3	A12	567	77056039	G	A ^b	9.6%	8.4×10 ⁻¹¹	
FS4	D04	23	53480934	C	T ^b	6.8%	2.0×10 ⁻⁹	
FS5	D06	1	5570632	T	C ^b	17.1%	1.6×10 ⁻⁷	Ma <i>et al.</i> , 2018
FS6	D07	2	39082751	T	G ^b	7.6%	4.6×10 ⁻⁸	Huang <i>et al.</i> , 2017
FE1	A05	1	38430737	T ^b	C	6.8%	1.3×10 ⁻⁸	
FE2	A07	4	93524441	C ^b	T	8.4%	4.1×10 ⁻⁸	Ma <i>et al.</i> , 2018
FE3	A08	1	92006893	T ^b	A	5.2%	6.3×10 ⁻¹¹	
FE4	A12	16	91652556	A	G ^b	29.1%	2.9×10 ⁻⁸	
FE5	A13	6	5588715	C	T ^b	24.7%	1.6×10 ⁻⁷	
FE6	D02	1	61894966	G	A ^b	8.0%	5.9×10 ⁻⁸	
FE7	D04	64	52570413	C	T ^b	48.2%	8.9×10 ⁻¹¹	Huang <i>et al.</i> , 2017; Thyssen <i>et al.</i> , 2019; Wang <i>et al.</i> , 2017
FE8	D05	71	27233958	T ^b	C	5.6%	2.3×10 ⁻⁸	Thyssen <i>et al.</i> , 2019
FU1	A01	2	2247417	G ^b	A	13.1%	7.1×10 ⁻⁸	
FU2	A09	5	81113800	G ^b	T	8.0%	7.2×10 ⁻¹⁰	
FU3	D05	28	27260152	G ^b	A	5.6%	1.2×10 ⁻⁹	Ma <i>et al.</i> , 2018
FU4	D12	21	46644449	G ^b	A	33.1%	3.9×10 ⁻⁸	

^aThe number of significant single nucleotide polymorphism (SNP).

^bFavorable SNP alleles.

^cOverlapping GWAS loci with previous studies.

Table 2 Identification of significant gene-trait associations in *Gossypium hirsutum* using transcriptome-wide association study (TWAS).

QTL	GWAS lead SNP position	GWAS <i>P</i> -value	Significant TWAS gene	TWAS FDR/	Homolog	Description ^a
FL6	D01:9210491	2.8×10 ⁻⁸	Ghir_D01G006540	9.2×10 ⁻³	AT4G27270	Quinone reductase family protein
FL7	D05:27233958	5.2×10 ⁻⁹	Ghir_D05G027990	2.3×10 ⁻⁴	AT1G22640	MYB domain protein 3
FL9	D11:24616418	6.6×10 ⁻¹⁰	Ghir_D11G020340	2.9×10 ⁻²	AT3G19150	KIP-related protein 6
			Ghir_D11G020430	2.4×10 ⁻⁴	AT1G09760	U2 small nuclear ribonucleoprotein A
FE5	A13:5588715	1.6×10 ⁻⁷	Ghir_A13G004410	5.5×10 ⁻⁴	AT2G14820	Phototropic-responsive NPH3 family protein
FE7	D04:52570413	8.9×10 ⁻¹¹	Ghir_D04G017090	1.3×10 ⁻⁶	AT4G36720	HVA22-like protein K
			Ghir_D04G016300	5.9×10 ⁻³	AT1G21070	Nucleotide-sugar transporter family protein
FU3	D05:27260152	1.2×10 ⁻⁹	Ghir_D05G027990	7.4×10 ⁻⁴	AT1G22640	MYB domain protein 3
FU4	D12:46644449	3.9×10 ⁻⁸	Ghir_D12G015450	6.0×10 ⁻⁴	AT5G50850	Transketolase family protein
			Ghir_D12G015670	6.3×10 ⁻³	AT3G48000	Aldehyde dehydrogenase 2B4

^aThe characterized function of homologous genes in *Arabidopsis*.

Figure legends

Figure 1 Genome-wide association study (GWAS) on cotton fiber quality-related traits.

These traits include fiber length (FL), fiber strength (FS), fiber elongation (FE) and fiber uniformity (FU). This analysis was performed using 2,372,767 single nucleotide polymorphisms (SNPs) in 251 cotton accessions. The horizontal red lines show the significance threshold of GWAS ($1/n$; 6.4). The x-axes show the 26 chromosomes (A01–A13 and D01–D13) in *Gossypium hirsutum*. Each chromosome is scaled by the physical chromosome length. The significant GWAS associations are also shown in **Table 1**.

Figure 2 Identification of eQTL using RNA-Seq data in cotton fiber development.

(a) Dot-plot showing eQTL and their regulated genes in 26 chromosomes. X-axis shows the single nucleotide polymorphism (SNP) position (bp) in each chromosome and y-axis shows gene position (bp) in each chromosome, with a chromosome order of A01 to D13 from left to right (x-axis) or from lower to upper (y-axis). The color of each dot represents the significance (P -value) of each eQTL-gene association, with low significance in green and high significance in blue. Each chromosome is scaled by the physical chromosome length. Dots in the diagonal line show the intra-chromosomal associations. The three red arrows show the enrichment of inter-subgenomic associations in homoeologous chromosomes and inter-chromosomal associations between the D11 and other chromosomes. **(b)** The difference of explanation rate (r^2) of SNPs for expression variation between local eQTL (<1 Mb) and distal eQTL (>1 Mb). The violin plots show the distribution density and box plots show the distribution quantiles. Two-sided Wilcoxon rank sum test, $**P$ -value $< 2.2 \times 10^{-16}$. **(c)** The distribution of distance (<100 Kb) between eQTL and regulated genes. **(d)** The proportions of local eQTL and distal eQTL. The distal eQTL were divided into three groups, including associations between the A-subgenome (At) and D-subgenome (Dt; At–Dt), between the At and At (At–At) and between the Dt and Dt (Dt–Dt). **(e)** The distribution of the number of eQTL for genes which were regulated by eQTL.

Figure 3 Analysis of eQTL and their regulated genes at the subgenome level in *Gossypium hirsutum*.

(a) Summary of eQTL in the A- and D-subgenomes involved in intra-subgenomic or inter-subgenomic regulation. (b) Summary of eGenes in the A- and D-subgenomes which were regulated by intra-subgenomic or inter-subgenomic eQTL. For a and b, the percentages show the proportions of intra-/inter-subgenomic eQTL in all distal eQTL or eGenes regulated by intra-/inter-subgenomic eQTL. (c) GO enrichment of genes in the A-subgenome (At-eGenes) which were regulated by eQTL in the D-subgenome (yellow bar charts) and genes in the D-subgenome (Dt-eGenes) which were regulated by eQTL in the A-subgenome (green bar charts). For each subgenome, all expressed genes at this fiber developmental stage were used as a reference for GO enrichment analysis.

Figure 4 Construction of eQTL regulatory network involving of eQTL hotspots.

(a) Genetic network between eQTL and genes. The green circle nodes represent genes which are regulated by eQTL, the yellow triangle nodes represent eQTL, and the octagon nodes represent eQTL hotspots. The eQTL hotspot 216 (Hot216) on chromosome D11 is enlarged with a chromosome location from 24,432,352 bp to 24,627,170 bp. The blue network edges represent local eQTL associations and grey edges represent distal eQTL associations. For Hot216, only distal eQTL associations are shown. (b) Summary of the number of Hot216-involved associations in different categories. The three different symbols have the same meanings as those in panel a. (c) Gene ontology (GO) enrichment of genes which are regulated by Hot216. These GO terms include those representing cellular component (CC), molecular function (MF) and biological processes (BP).

Figure 5 Characterization of eQTL hotspot 216 (Hot216) on chromosome D11.

(a) Clustering analysis of the expression levels of eGenes which are regulated by Hot216. The 962 genes are categorized into two groups (group-I with 287 genes and group-II with 675 genes) and the 251 cotton accessions are clustered into two clusters (cluster-I with 91 accessions and cluster-II with 160 accessions). The significant GO terms in each gene group are shown. (b) Genome-wide distribution of eGenes regulated by Hot216. Some representative genes are shown with the full description of gene names in Table S8. (c) Manhattan plot of GWAS signal on chromosome D11 which is associated with fiber length. (d) Manhattan plot of the eQTL signal of *Ghir_D11G020340*.

For **C** and **D**, the horizontal red lines show the significance threshold ($1/n$; 6.4). **(e)** Distribution of fiber length with two different genotypes of *Ghir_D11G020340*. Two-sided Wilcoxon rank sum test, $**P\text{-value} < 2.2 \times 10^{-16}$. The horizontal lines in box plots show median values, and ranges show the first and third quartiles. **(f)** Normalized expression (FPKM) of representative genes in two cotton groups categorized by the SNP site of *Ghir_D11G020340*.

Figure 6 Expression analysis of *CesA* genes in cotton accessions with extreme long or short fibers. 30 cotton accessions with extreme long fibers were compared with 30 accessions with extreme short fibers to identify differentially expressed genes. The long-fiber and short-fiber accessions are the same as those in **Fig. S4**. The horizontal lines in box plots show median values, and ranges show the first and third quartiles. All these *CesA* genes show differential expression levels between the two groups of accessions ($**\text{FDR} \leq 0.05$).

References

- Adams, K.L., Cronn, R., Percifield, R., and Wendel, J.F. 2003. Genes duplicated by polyploidy show unequal contributions to the transcriptome and organ-specific reciprocal silencing. *Proc. Natl Acad. Sci. USA* **100**:4649-4654.
- Aguet, F., Brown, A.A., Castel, S.E., Davis, J.R., H, Y., J, B., Mohammadi, P., Park, Y.S., Parsana, P., Segrè, A.V., *et al.* 2017. Genetic effects on gene expression across human tissues. *Nature* **550**:204-213.
- Anders, S., and Huber, W. 2010. Differential expression analysis for sequence count data. *Genome Biol.* **11**:R106.
- Barrôco, R.M., Peres, A., Droual, A.-M., De Veylder, L., Nguyen, L.S.L., De Wolf, J., Mironov, V., Peerbolte, R., Beemster, G.T.S., Inzé, D., *et al.* 2006. The Cyclin-Dependent Kinase Inhibitor Orysa;KRP1 Plays an Important Role in Seed Development of Rice. *Plant Physiol.* **142**:1053-1064.
- Bischoff, V., Nita, S., Neumetzler, L., Schindelasch, D., Urbain, A., Eshed, R., Persson, S., Delmer, D., and Scheible, W.R. 2010. *TRICHOME BIREFRINGENCE* and its homolog *AT5G01360* encode plant-specific DUF231 proteins required for cellulose biosynthesis in *Arabidopsis*. *Plant Physiol.* **153**:590-602.
- Bolger, A.M., Lohse, M., and Usadel, B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**:2114-2120.
- Bourquin, V., Nishikubo, N., Abe, H., Brumer, H., Denman, S., Eklund, M., Christiernin, M., Teeri, T.T., Sundberg, B., and Mellerowicz, E.J. 2002. Xyloglucan Endotransglycosylases Have a Function during the Formation of Secondary Cell Walls of Vascular Tissues. *Plant Cell* **14**:3073-3088.
- Cheung, V.G., and Spielman, R.S. 2009. Genetics of human gene expression: mapping DNA variants that influence gene expression. *Nat. Rev. Genet.* **10**:595-604.
- Cookson, W., Liang, L., Abecasis, G., Moffatt, M., and Lathrop, M. 2009. Mapping complex disease traits with global gene expression. *Nat. Rev. Genet.* **10**:184-194.

- Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker, R.E., Lunter, G., Marth, G.T., Sherry, S.T., *et al.* 2011. The variant call format and VCFtools. *Bioinformatics* **27**:2156-2158.
- Falush, D., Stephens, M. & Pritchard, J.K. 2003. Inference of population structure using multilocus genotype data: Linked loci and correlated allele frequencies. *Genetics* **164**:1567–1587.
- Fang, L., Wang, Q., Hu, Y., Jia, Y., Chen, J., Liu, B., Zhang, Z., Guan, X., Chen, S., Zhou, B., *et al.* 2017. Genomic analyses in cotton identify signatures of selection and loci associated with fiber quality and yield traits. *Nat. Genet.* **49**:1089-1098.
- Fu, J.J., Cheng, Y.B., Linghu, J.J., Yang, X.H., Kang, L., Zhang, Z.X., Zhang, J., He, C., Du, X.M., Peng, Z.Y., *et al.* 2013. RNA sequencing reveals the complex regulatory network in the maize kernel. *Nat. Commun.* **4**:2832.
- García-Ramírez, E., Rosas-Bringas, O.G., Godínez-Palma, S.K., Vázquez-Ramos, J.M., Rosas-Bringas, F.R., and Zamora-Zaragoza, J. 2017. Two maize Kip-related proteins differentially interact with, inhibit and are phosphorylated by cyclin D–cyclin-dependent kinase complexes. *J. Exp. Bot.* **68**:1585-1597.
- Guo, K., Du, X., Tu, L., Tang, W., Wang, P., Wang, M., Liu, Z., and Zhang, X. 2016. Fibre elongation requires normal redox homeostasis modulated by cytosolic ascorbate peroxidase in cotton (*Gossypium hirsutum*). *J. Exp. Bot.* **67**:3289-3301.
- Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B.W., Jansen, R., de Geus, E.J., Boomsma, D.I., Wright, F.A., *et al.* 2016. Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.* **48**:245-252.
- Gusev, A., Mancuso, N., Won, H., Kousi, M., Finucane, H.K., Reshef, Y., Song, L., Safi, A., Schizophrenia Working Group of the Psychiatric Genomics, C., McCarroll, S., *et al.* 2018. Transcriptome-wide association study of schizophrenia and chromatin activity yields mechanistic disease insights. *Nat. Genet.* **50**:538-548.
- Haigler, C., Betancur, L., Stiff, M., and Tuttle, J. 2012. Cotton fiber: a powerful single-cell model for cell wall and cellulose research. *Front Plant Sci.* **3**:104.
- Hakonarson, H., Li, M., and Wang, K. 2010. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**:e164-e164.

- Han, L.B., Li, Y.B., Wang, H.Y., Wu, X.M., Li, C.L., Luo, M., Wu, S.J., Kong, Z.S., Pei, Y., Jiao, G.L., et al.** 2013. The Dual Functions of WLIM1a in Cell Elongation and Secondary Wall Formation in Developing Cotton Fibers. *Plant Cell* **25**:4421-4438.
- Holanda, A.J., Silva, I.T., Nussenzweig, M.C., Jankovic, M., and Rosales, R.A.** 2014. Identification of chromosomal translocation hotspots via scan statistics. *Bioinformatics* **30**:2551-2558.
- Hormozdiari, F., Gazal, S., van de Geijn, B., Finucane, H.K., Ju, C.J.T., Loh, P.-R., Schoech, A., Reshef, Y., Liu, X., O'Connor, L., et al.** 2018. Leveraging molecular quantitative trait loci to understand the genetic architecture of diseases and complex traits. *Nat. Genet.* **50**:1041-1047.
- Hovav, R., Udall, J.A., Chaudhary, B., Rapp, R., Flagel, L., and Wendel, J.F.** 2008. Partitioned expression of duplicated genes during development and evolution of a single cell in a polyploid plant. *Proc. Natl Acad. Sci. USA* **105**:6191-6195.
- Huang, C., Nie, X., Shen, C., You, C., Li, W., Zhao, W., Zhang, X., and Lin, Z.** 2017. Population structure and genetic basis of the agronomic traits of upland cotton in China revealed by a genome-wide association study using high-density SNPs. *Plant Biotechnol J* **15**:1374-1386.
- Jégu, T., Latrasse, D., Delarue, M., Mazubert, C., Bourge, M., Hudik, E., Blanchet, S., Soler, M.-N., Charon, C., De Veylder, L., et al.** 2013. Multiple Functions of Kip-Related Protein5 Connect Endoreduplication and Cell Elongation. *Plant Physiol* **161**:1694-1705.
- Keurentjes, J.J.B., Fu, J., Terpstra, I.R., Garcia, J.M., van den Ackerveken, G., Snoek, L.B., Peeters, A.J.M., Vreugdenhil, D., Koornneef, M., and Jansen, R.C.** 2007. Regulatory network construction in *Arabidopsis* by using genome-wide gene expression quantitative trait loci. *Proc. Natl Acad. Sci. USA* **104**:1708-1713.
- Kim, D., Langmead, B., and Salzberg, S.L.** 2015. HISAT: a fast spliced aligner with low memory requirements. *Nat. Methods* **12**:357-360.
- Kim, W.-C., Ko, J.-H., Kim, J.-Y., Kim, J., Bae, H.-J., and Han, K.-H.** 2013. MYB46 directly regulates the gene expression of secondary wall-associated cellulose synthases in *Arabidopsis*. *Plant J.* **73**:26-36.

- Kremling, K.A.G., Chen, S.-Y., Su, M.-H., Lepak, N.K., Romay, M.C., Swarts, K.L., Lu, F., Lorant, A., Bradbury, P.J., and Buckler, E.S. 2018. Dysregulation of expression correlates with rare-allele burden and fitness loss in maize. *Nature* **555**:520-523.
- Kroon, D.E., Buckler, E.S., Bradbury, P.J., Casstevens, T.M., Ramdoss, Y., and Zhang, Z. 2007. TASSEL: software for association mapping of complex traits in diverse samples. *Bioinformatics* **23**:2633-2635.
- Lappalainen, T., Sammeth, M., Friedländer, M.R., 't Hoen, P.A.C., Monlong, J., Rivas, M.A., González-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G., *et al.* 2013. Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501**:506-511.
- Lee, C., O'Neill, M.A., Tsumuraya, Y., Darvill, A.G., and Ye, Z.-H. 2007. The irregular xylem9 Mutant is Deficient in Xylan Xylosyltransferase Activity. *Plant Cell Physiol.* **48**:1624-1634.
- Li, H., and Durbin, R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**:1754-1760.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**:2078-2079.
- Li, Y., Tu, L., Pettolino, F.A., Ji, S., Hao, J., Yuan, D., Deng, F., Tan, J., Hu, H., Wang, Q., *et al.* 2016. GbEXPATR, a species-specific expansin, enhances cotton fibre elongation through cell wall restructuring. *Plant Biotechnol. J.* **14**:951-963.
- Lippert, C., Listgarten, J., Liu, Y., Kadie, C.M., Davidson, R.I., and Heckerman, D. 2011. FaST linear mixed models for genome-wide association studies. *Nat. Methods* **8**:833-835.
- Liu, H., Luo, X., Niu, L., Xiao, Y., Chen, L., Liu, J., Wang, X., Jin, M., Li, W., Zhang, Q., *et al.* 2017. Distant eQTLs and non-coding sequences play critical roles in regulating gene expression and quantitative trait variation in maize. *Mol. Plant* **10**:414-426.
- Liu, R., Gong, J., Xiao, X., Zhang, Z., Li, J., Liu, A., Lu, Q., Shang, H., Shi, Y., Ge, Q., *et al.* 2018. GWAS Analysis and QTL Identification of Fiber Quality Traits and Yield Components in Upland Cotton Using Enriched High-Density SNP Markers. *Front Plant Sci.* **9**:1067-1067.
- Lowry, D.B., Logan, T.L., Santuari, L., Hardtke, C.S., Richards, J.H., DeRose-Wilson, L.J., McKay, J.K., Sen, S., and Juenger, T.E. 2013. Expression Quantitative Trait Locus

Mapping across Water Availability Environments Reveals Contrasting Associations with Genomic Features in *Arabidopsis*. *Plant Cell* **25**:3266-3279.

Ma, Z., He, S., Wang, X., Sun, J., Zhang, Y., Zhang, G., Wu, L., Li, Z., Liu, Z., Sun, G., et al. 2018. Resequencing a core collection of upland cotton identifies genomic variation and loci influencing fiber quality and yield. *Nat. Genet.* **50**:803-813.

MacMillan, C.P., Mansfield, S.D., Stachurski, Z.H., Evans, R., and Southerton, S.G. 2010. Fasciclin-like arabinogalactan proteins: specialization for stem biomechanics and cell wall architecture in *Arabidopsis* and *Eucalyptus*. *Plant J.* **62**:689-703.

McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M., et al. 2010. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**:1297-1303.

Pertea, M., Pertea, G.M., Antonescu, C.M., Chang, T.-C., Mendell, J.T., and Salzberg, S.L. 2015. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.* **33**:290-295.

Sablowski, R., and Carnier Dornelas, M. 2013. Interplay between cell growth and cell cycle in plants. *J. Exp. Bot.* **65**:2703-2714.

Shan, C.-M., Shangguan, X.-X., Zhao, B., Zhang, X.-F., Chao, L.-m., Yang, C.-Q., Wang, L.-J., Zhu, H.-Y., Zeng, Y.-D., Guo, W.-Z., et al. 2014. Control of cotton fibre elongation by a homeodomain transcription factor GhHOX3. *Nat. Commun.* **5**:5519.

Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., and Ideker, T. 2003. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**:2498-2504.

Tan, J., Tu, L., Deng, F., Hu, H., Nie, Y., and Zhang, X. 2013. A Genetic and Metabolic Analysis Revealed that Cotton Fiber Cell Development Was Retarded by Flavonoid Naringenin. *Plant Physiol.* **162**:86-95.

Taylor-Teeples, M., Lin, L., de Lucas, M., Turco, G., Toal, T.W., Gaudinier, A., Young, N.F., Trabucco, G.M., Veling, M.T., Lamothe, R., et al. 2015. An *Arabidopsis* gene regulatory network for secondary cell wall synthesis. *Nature* **517**:571-575.

- Thyssen, G.N., Jenkins, J.N., McCarty, J.C., Zeng, L., Campbell, B.T., Delhom, C.D., Islam, M.S., Li, P., Jones, D.C., Condon, B.D., *et al.* 2019. Whole genome sequencing of a MAGIC population identified genomic loci and candidate genes for major fiber quality traits in upland cotton (*Gossypium hirsutum* L.). *Theor Appl Genet.* 132:989-999.
- Tian, T., Liu, Y., Yan, H., You, Q., Yi, X., Du, Z., Xu, W., and Su, Z. 2017. agriGO v2.0: a GO analysis toolkit for the agricultural community, 2017 update. *Nucleic Acids Res.* 45:W122-W129.
- Tuttle, J.R., Nah, G., Duke, M.V., Alexander, D.C., Guan, X., Song, Q., Chen, Z.J., Scheffler, B.E., and Haigler, C.H. 2015. Metabolomic and transcriptomic insights into how cotton fiber transitions to secondary wall synthesis, represses lignification, and prolongs elongation. *BMC Genomics* 16, 477.
- Van't Hof, J. 1999. Increased nuclear DNA content in developing cotton fiber cells. *Am. J. Bot.* 86:776-779.
- Wang, J., Yu, H., Weng, X., Xie, W., Xu, C., Li, X., Xiao, J., and Zhang, Q. 2014. An expression quantitative trait loci-guided co-expression analysis for constructing regulatory network using a rice recombinant inbred line population. *J. Exp. Bot.* 65:1069-1079.
- Wang, J., Yu, H., Xie, W., Xing, Y., Yu, S., Xu, C., Li, X., Xiao, J., and Zhang, Q. 2010. A global analysis of QTLs for expression variations in rice shoots at the early seedling stage. *Plant J.* 63:1063-1074.
- Wang, M., Tu, L., Lin, M., Lin, Z., Wang, P., Yang, Q., Ye, Z., Shen, C., Li, J., Zhang, L., *et al.* 2017. Asymmetric subgenome selection and *cis*-regulatory divergence during cotton domestication. *Nat. Genet.* 49:579-587.
- Wang, M., Tu, L., Yuan, D., Zhu, D., Shen, C., Li, J., Liu, F., Pei, L., Wang, P., Zhao, G., *et al.* 2019. Reference genome sequences of two cultivated allotetraploid cottons, *Gossypium hirsutum* and *Gossypium barbadense*. *Nat. Genet.* 51:224-229.
- Wang, M., Wang, P., Lin, M., Ye, Z., Li, G., Tu, L., Shen, C., Li, J., Yang, Q., and Zhang, X. 2018. Evolutionary dynamics of 3D genome architecture following polyploidization in cotton. *Nat. Plants* 4:90-97.

- Wang, X., Chen, Q., Wu, Y., Lemmon, Z.H., Xu, G., Huang, C., Liang, Y., Xu, D., Li, D., Doebley, J.F., *et al.* 2018. Genome-wide Analysis of Transcriptional Variability in a Large Maize-Teosinte Population. *Mol. Plant* **11**:443-459.
- West, M.A.L., Kim, K., Kliebenstein, D.J., van Leeuwen, H., Michelmore, R.W., Doerge, R.W., and St Clair, D.A. 2007. Global eQTL mapping reveals the complex genetic architecture of transcript-level variation in *Arabidopsis*. *Genetics* **175**:1441-1450.
- Yang, J., Lee, S.H., Goddard, M.E., and Visscher, P.M. 2011. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**:76-82.
- Yuan, D., Tang, Z., Wang, M., Gao, W., Tu, L., Jin, X., Chen, L., He, Y., Zhang, L., Zhu, L., *et al.* 2015. The genome sequence of Sea-Island cotton (*Gossypium barbadense*) provides insights into the allopolyploidization and development of superior spinnable fibres. *Sci. Rep.* **5**:17662.
- Zhang, L., Su, W., Tao, R., Zhang, W., Chen, J., Wu, P., Yan, C., Jia, Y., Larkin, R.M., Lavelle, D., *et al.* 2017. RNA sequencing provides insights into the evolution of lettuce and the regulation of flavonoid biosynthesis. *Nat. Commun.* **8**:2264.
- Zhang, T., Hu, Y., Jiang, W., Fang, L., Guan, X., Chen, J., Zhang, J., Saski, C.A., Scheffler, B.E., *et al.* 2015. Sequencing of allotetraploid cotton (*Gossypium hirsutum* L. acc. TM-1) provides a resource for fiber improvement. *Nat Biotechnol* **33**:531-537.
- Zhang, X., Cal, A.J., and Borevitz, J.O. 2011. Genetic architecture of regulatory variation in *Arabidopsis thaliana*. *Genome Res.* **21**:725-733.
- Zhang, Z., Ruan, Y.-L., Zhou, N., Wang, F., Guan, X., Fang, L., Shang, X., Guo, W., Zhu, S., and Zhang, T. 2017. Suppressing a Putative Sterol Carrier Gene Reduces Plasmodesmal Permeability and Activates Sucrose Transporter Genes during Cotton Fiber Elongation. *Plant Cell* **29**:2027-2046.
- Zhao, X.A., Bramsiepe, J., Van Durme, M., Komaki, S., Prusicki, M.A., Maruyama, D., Forner, J., Medzihradszky, A., Wijnker, E., Harashima, H., *et al.* 2017. RETINOBLASTOMA RELATED1 mediates germline entry in *Arabidopsis*. *Science* **356**:eaaf6532.
- Zhu, G., Wang, S., Huang, Z., Zhang, S., Liao, Q., Zhang, C., Lin, T., Qin, M., Peng, M., Yang, C., *et al.* 2018. Rewiring of the Fruit Metabolome in Tomato Breeding. *Cell* **172**:249-261.

Zhu, Z., Zhang, F., Hu, H., Bakshi, A., Robinson, M.R., Powell, J.E., Montgomery, G.W., Goddard, M.E., Wray, N.R., Visscher, P.M., *et al.* 2016. Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48**:481-487.

Supporting Information

Additional supporting information may be found in the online version of this article.

Figure S1 The comparison of significance (P-values) between intra-chromosomal eQTL (SameChr_eQTL) and inter-chromosomal eQTL (DiffChr_eQTL).

Figure S2 A visualization of GO hierarchical graphic design for three likely causal genes.

Figure S3 Comparison of fiber quality-related traits between cotton accessions in cluster1 and cluster2.

Figure S4 Differential expression analysis of genes in accessions with extreme fiber length.

Figure S5 Whole genome Hi-C matrix in *G. hirsutum*.

Table S1 Descriptive statistics for phenotypic variations and broad-sense heritability for 5 fiber quality-related traits.

Table S2 Summary of significant SNPs in GWAS loci.

Table S3 Summary of genes in GWAS loci with putative function in fiber development.

Table S4 Summary of RNA-Seq data in this study.

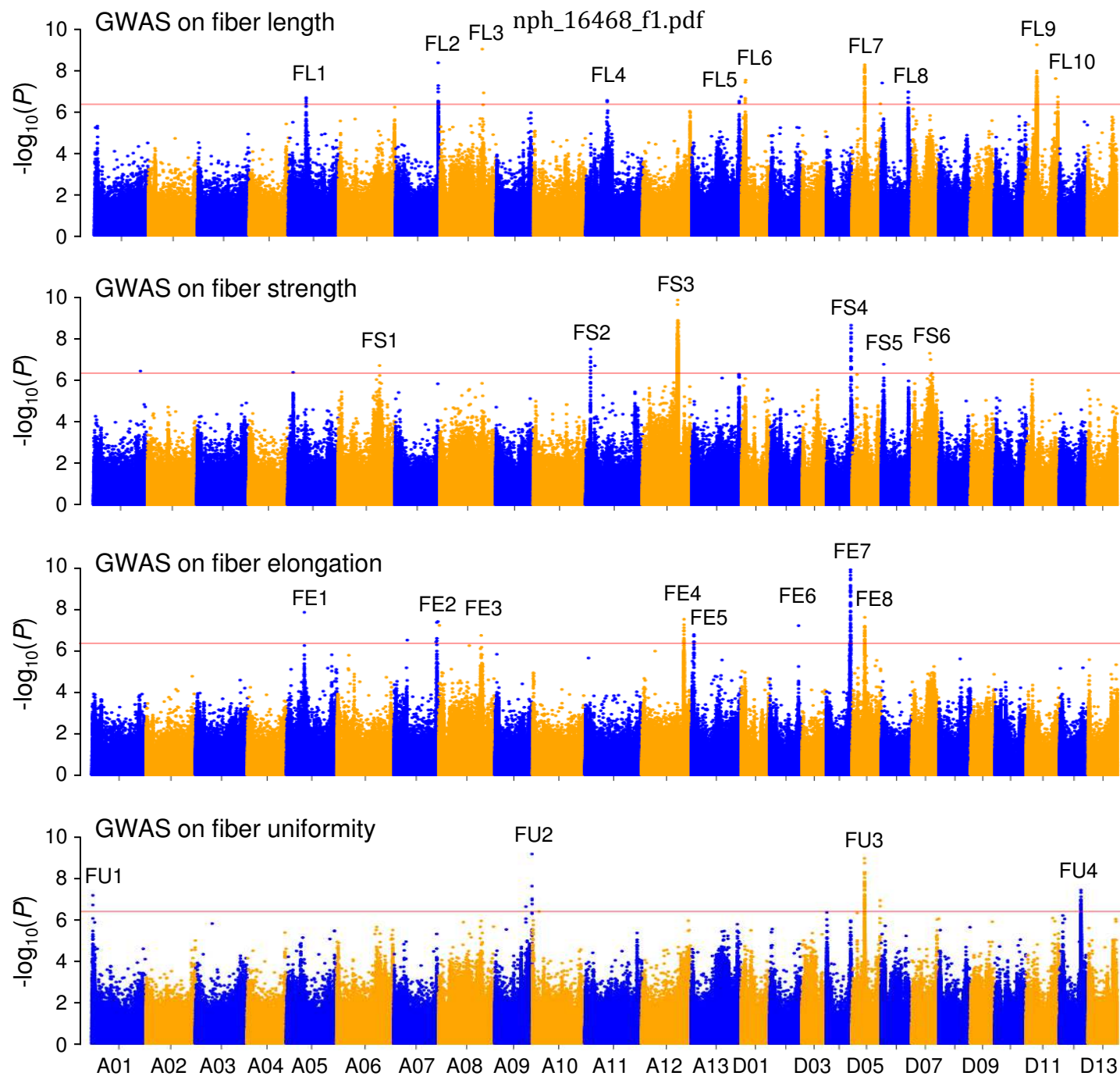
Table S5 GO enrichment of genes in one subgenome regulated by eQTL in the other subgenome.

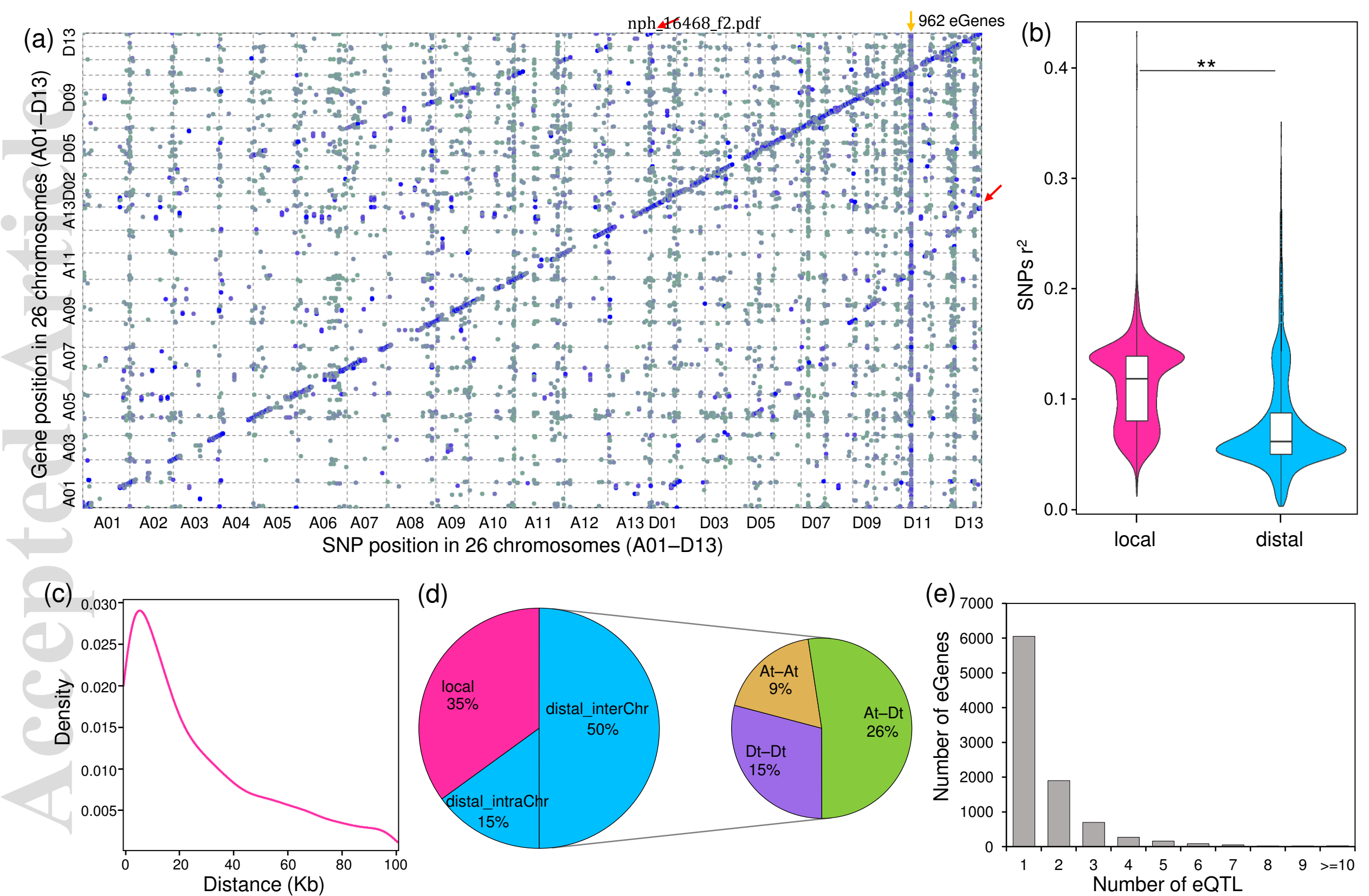
Table S6 Colocalization of GWAS associations and subgenomic eQTL regulation.

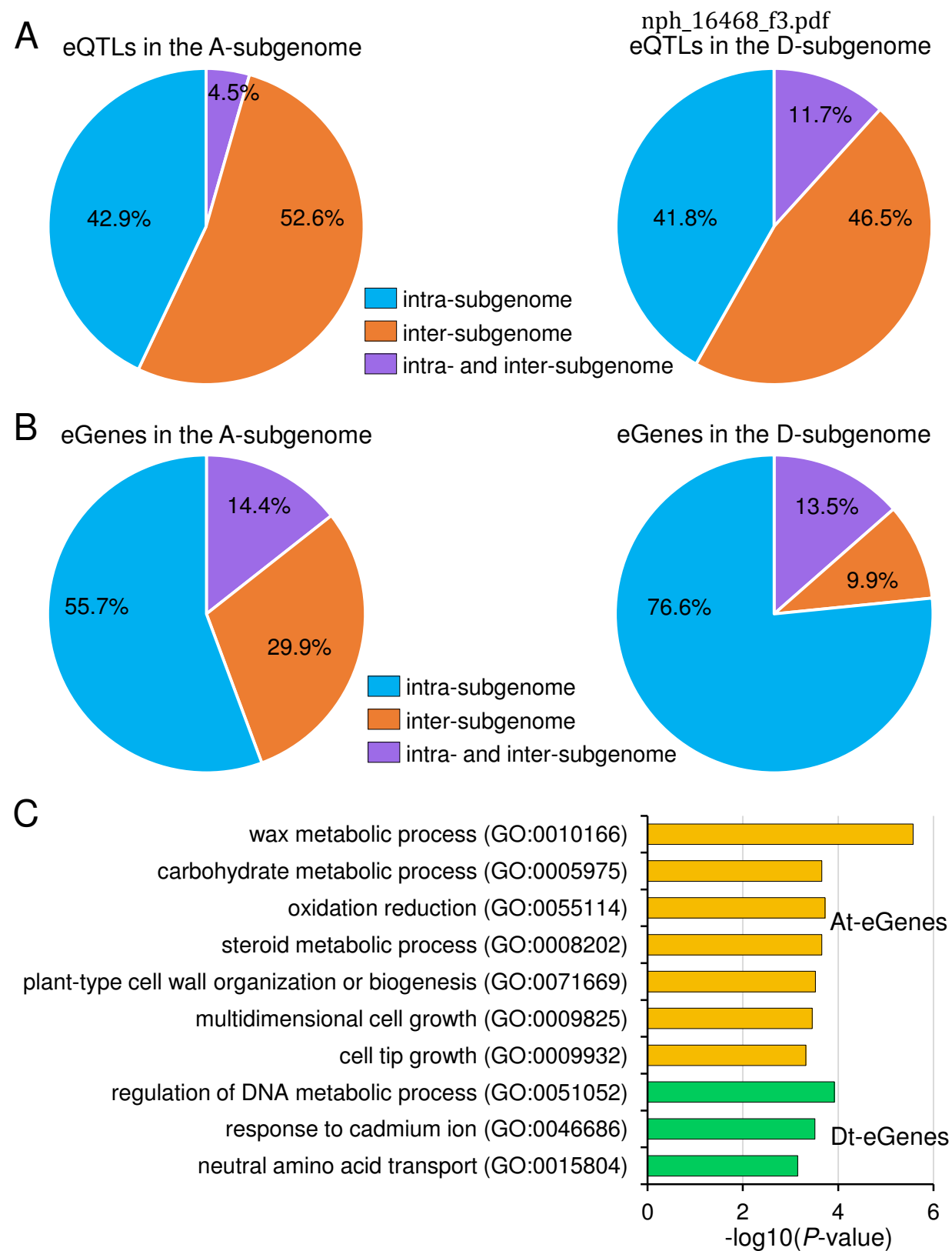
Table S7 Summary of eQTL hotspot in this study.

Table S8 Information for 962 eGenes regulated by the Hot216. Columns 2-6 show the expression levels (FPKM) of eGenes in fiber development.

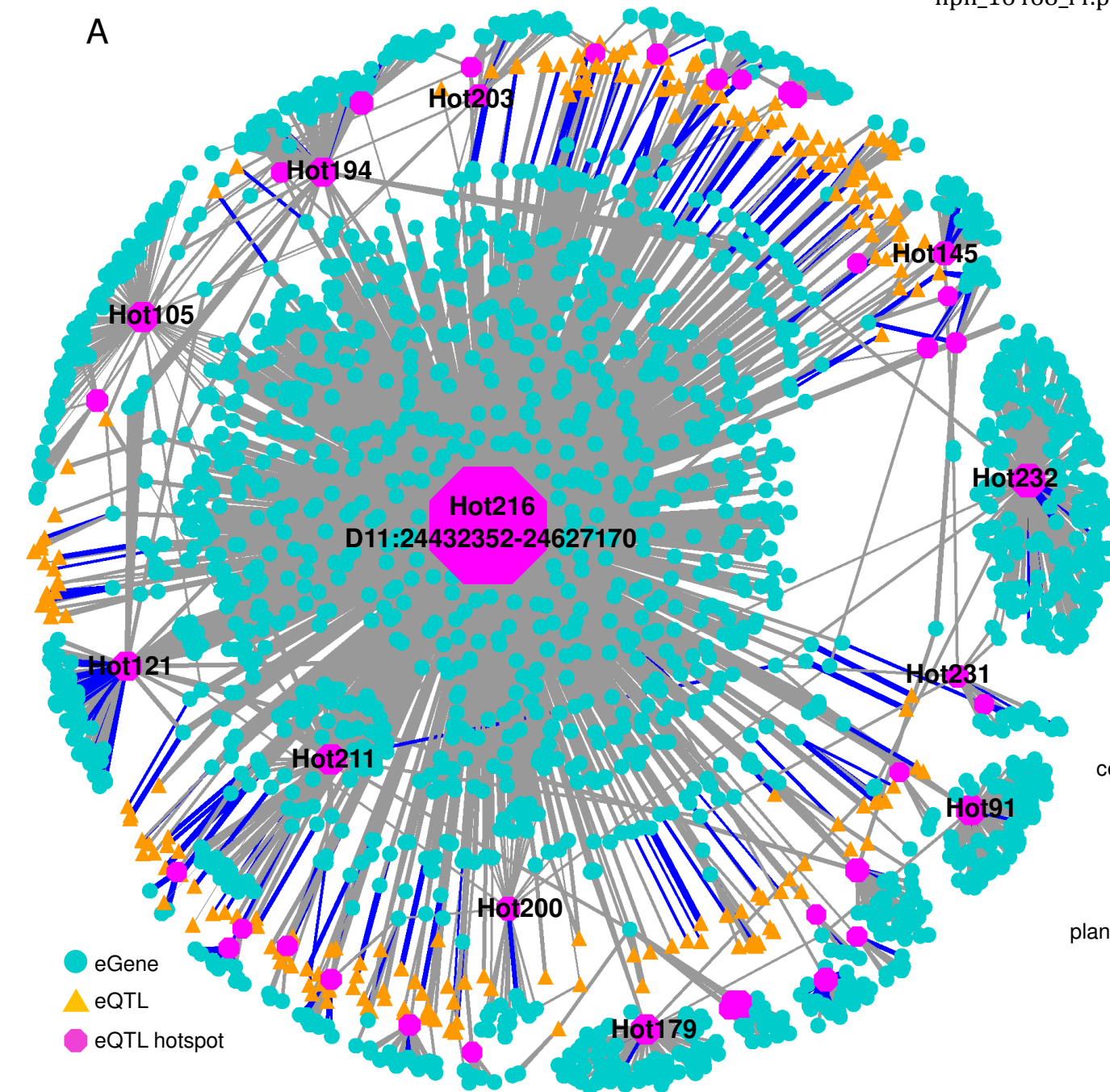
Table S9 GO enrichment of genes which were regulated by the eQTL hotspot in chromosome D11.



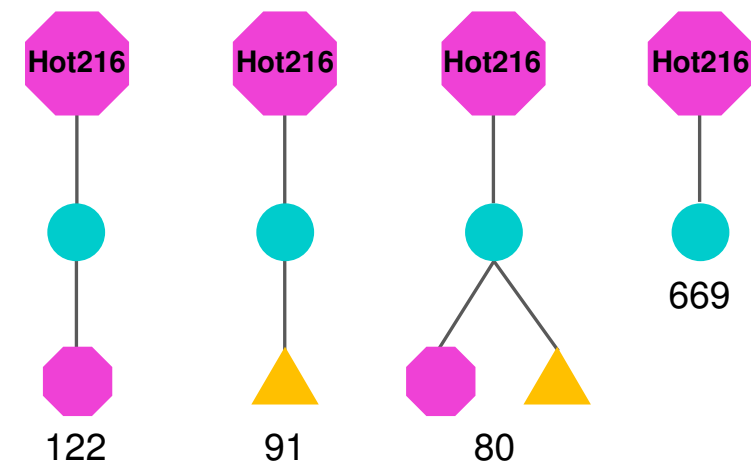




A



B



C

